# NSF OCI: #0940841 DataNet Federation Consortium
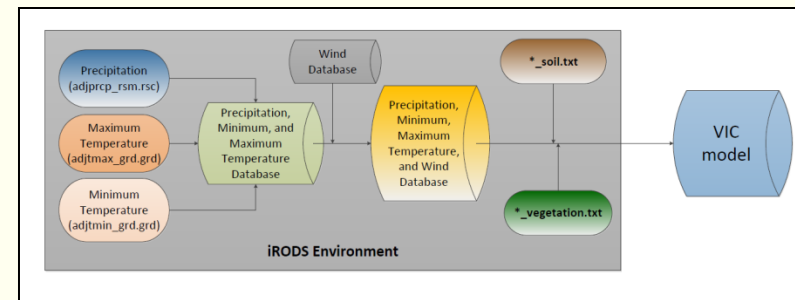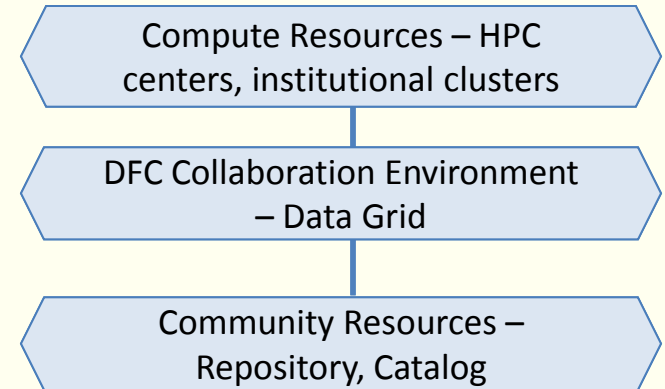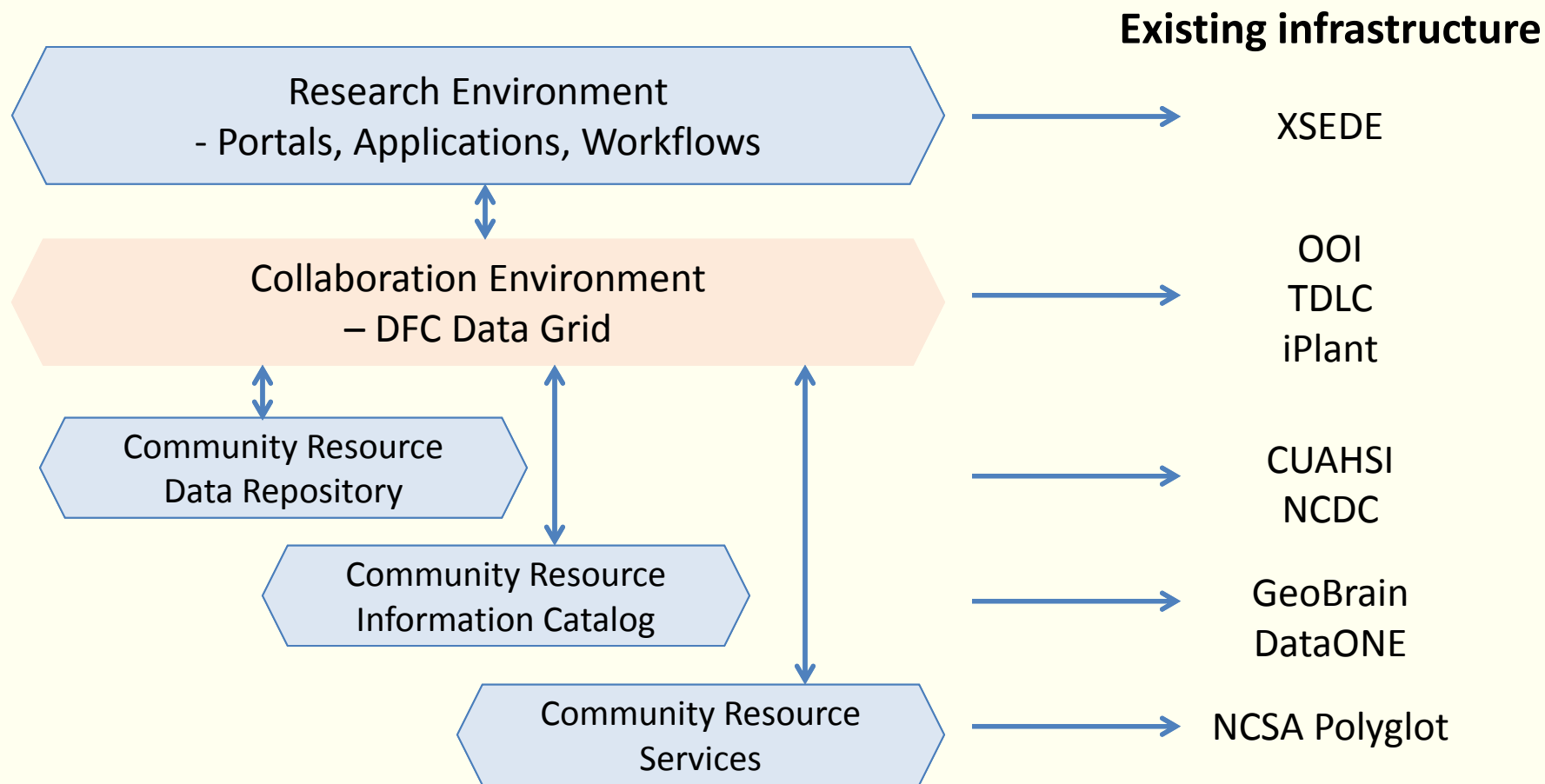
- **Enable collaborative research**
  - Sharing of data, information, and knowledge

- **Build national data cyberinfrastructure**
  - Federation of existing data management systems

- **Support reproducible data-driven research**
  - Encapsulate knowledge in shared workflows

- **Enable student participation in research**
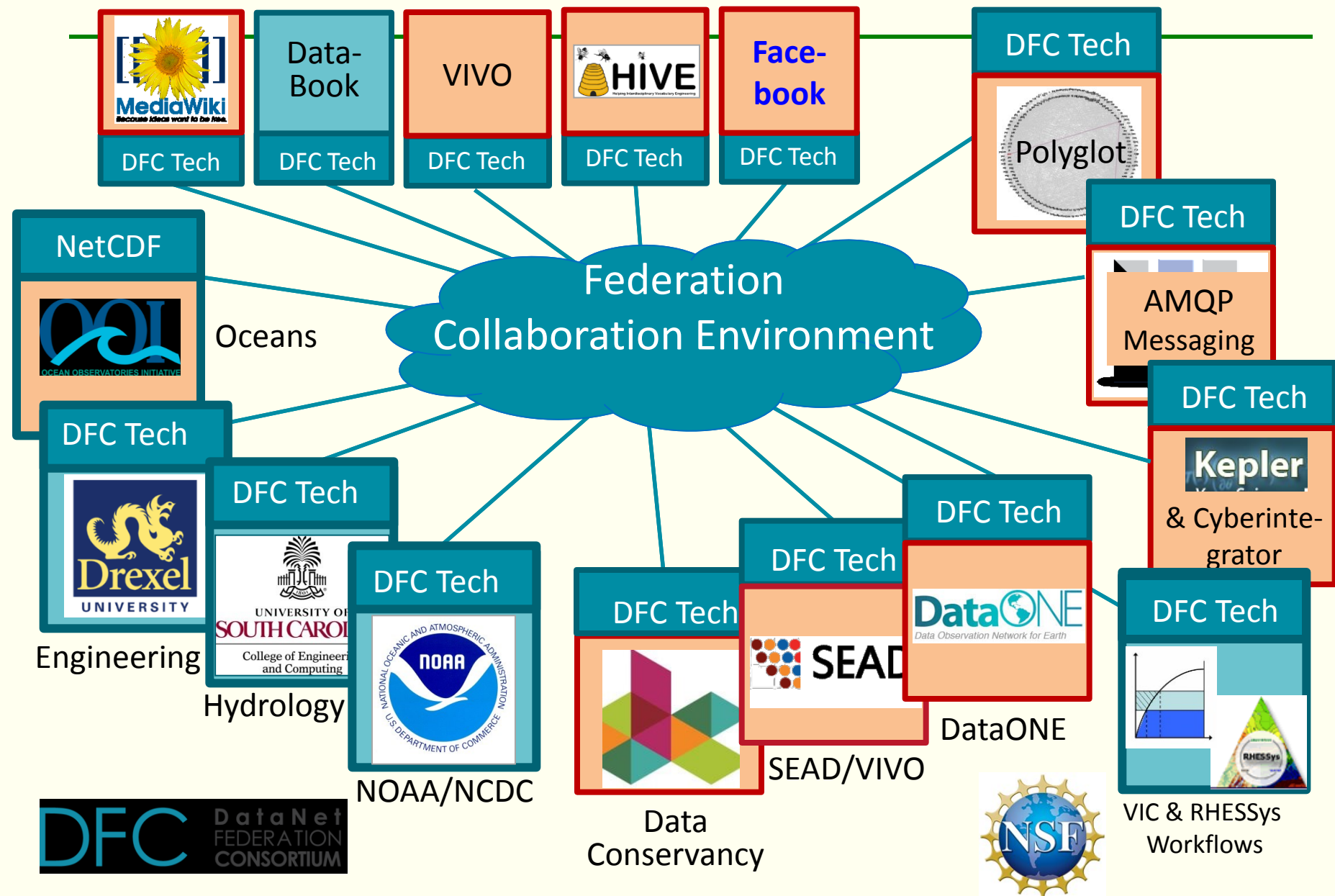  - Policy-controlled access to "live" data
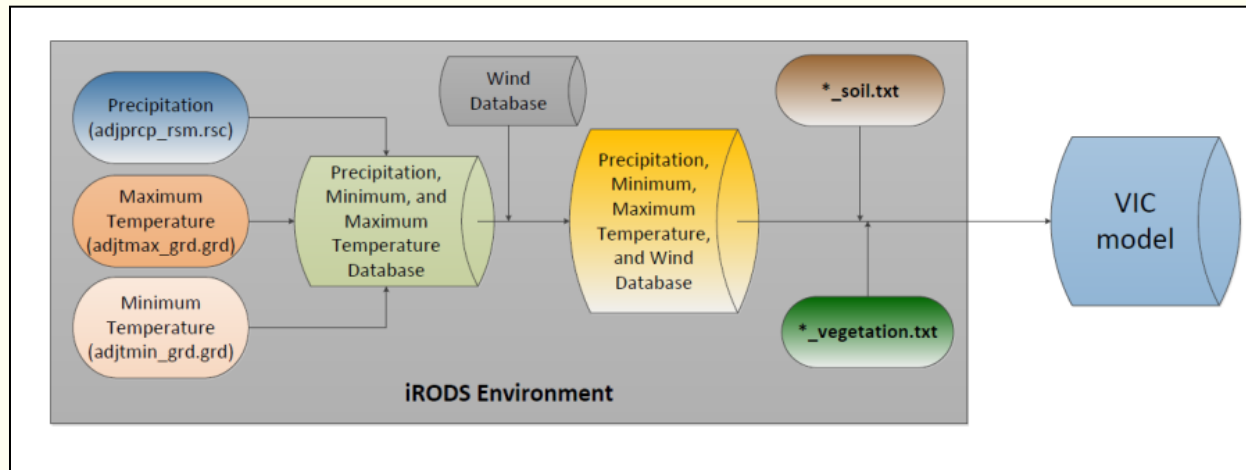
**NSF DataNet Federation Consortium**
*Enabling Collaboration through Interoperability*

DFC iRODS-based middleware enables interoperability between heterogeneous clients, data, and service resources

MediaWiki — DFC Tech
Data-Book — DFC Tech
VIVO — DFC Tech
HiVE — DFC Tech
Face-book — DFC Tech
DFC Tech — Polyglot
DFC Tech — AMQP Messaging

Federation Collaboration Environment

NetCDF
Oceans

DFC Tech — Drexel UNIVERSITY
Engineering

DFC Tech — UNIVERSITY OF SOUTH CAROLINA College of Engineering and Computing
Hydrology

DFC Tech — NOAA U.S DEPARTMENT OF COMMERCE
NOAA/NCDC

DFC Tech
Data Conservancy

DFC Tech — SEAD
SEAD/VIVO

DFC Tech — DataONE Data Observation Network for Earth
DataONE

DFC Tech — Kepler & Cyberinte-grator

DFC Tech — RHESSys
VIC & RHESSys Workflows

DFC DataNet FEDERATION CONSORTIUM

NSF

# Practitioners' Perspective

- Build community resource
  - Address explicit purpose for formation of a collaboration
  - Build community consensus on provenance, descriptive, system metadata
  - Capture domain knowledge (procedures for interoperability, research analyses, management)
  - Share data, procedures, workflows
- Enable reproducible data-driven research through workflows

# Challenges

- DFC uses iRODS policy-based data grid to handle:
  - Acquisition of all relevant data for research
    - Develop micro-services that can access external repositories
  - Distribution of data management effort
    - Use data grid to automate replication of data between agencies
  - Automation of the application of domain knowledge
    - Share workflows used in research analyses
  - Management of policies for data control
    - Enforce policies at each storage location

1. Metadata virtualization (manage properties of metadata – creation time, storage location, access controls, schema)

2. Knowledge virtualization (manage processes that generate metadata – provenance, descriptive, administrative)

# iRODS Policy-Based Data Management

- ***Purpose***               - reason a collection is assembled

- ***Properties***          - attributes needed to ensure the **purpose**

- ***Policies***              - rules to enforce and maintain collection **properties**

- ***Procedures***         - functions that implement the **policies**

- ***Persistent state information***     – metadata from applying the **procedures**

- ***Property assessment criteria***     – validation that **state information**

                           conforms to the desired **purpose**

- ***Federation***           - controlled sharing of **logical name spaces**


- We capture domain knowledge in policies and procedures, and evolve policies to implement data life cycle stages

- Broadening of impact corresponds to evolution of policies to represent consensus of a new larger community

# NSF Data Bridge: Solving the First & Last Mile Problems in Big Data

**First Mile:** Bring the Long-tail of Science Data into Mainstream

**Last Mile:** Automate Linking, Clustering, and Discovery of Interesting Relationships in Heterogeneous Data

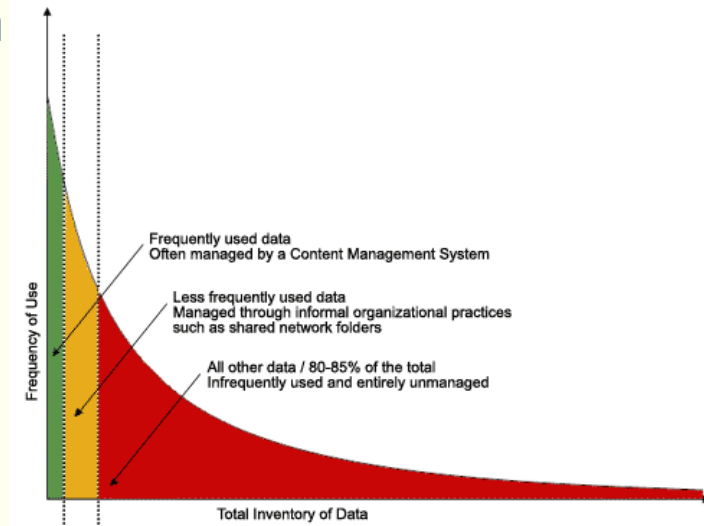**Data Bridge:** NSF-funded Big Data Project

– Apply Socio-metric Network Analysis (SNA) to data

– Explore Relationships between Data, Users, Resources, Methods, Workflows, …

– Link through Multi-dimensional vectors

- Similar to, but for data: **Linked in** · You Tube · f · Twitter · g+ · TAGGED

– Incentives:

- Enable participation in a larger collaboration
- Raise awareness of local data and bring low value per byte data into shared collections



Frequently used data
Often managed by a Content Management System

Less frequently used data
Managed through informal organizational practices such as shared network folders

All other data / 80–85% of the total
Infrequently used and entirely unmanaged

Frequency of Use

Total Inventory of Data

# More Information

- DataNet Federation Consortium
  - http://datafed.org
  - UNC-CH, UCSD, Drexel, USC
- Integrated Rule Oriented Data System (iRODS)
  - [http://irods.diceresearch.org](http://irods.diceresearch.org)
  - Application of data grids include
    - NOAA National Climatic Data Center
    - NASA Center for Climate Simulations
    - French National Library
    - Broad Institute genomics data grid
    - International Neuroinformatics Coordinating Facility